

IBIS2010

# 補助情報を用いた情報源符号化

松嶋敏泰  
早稲田大学

# 統計的学習理論と情報理論

(関係するものすぐ思いつくのは)

## 情報源符号化(冗長度を削減)

ユニバーサル情報源符号

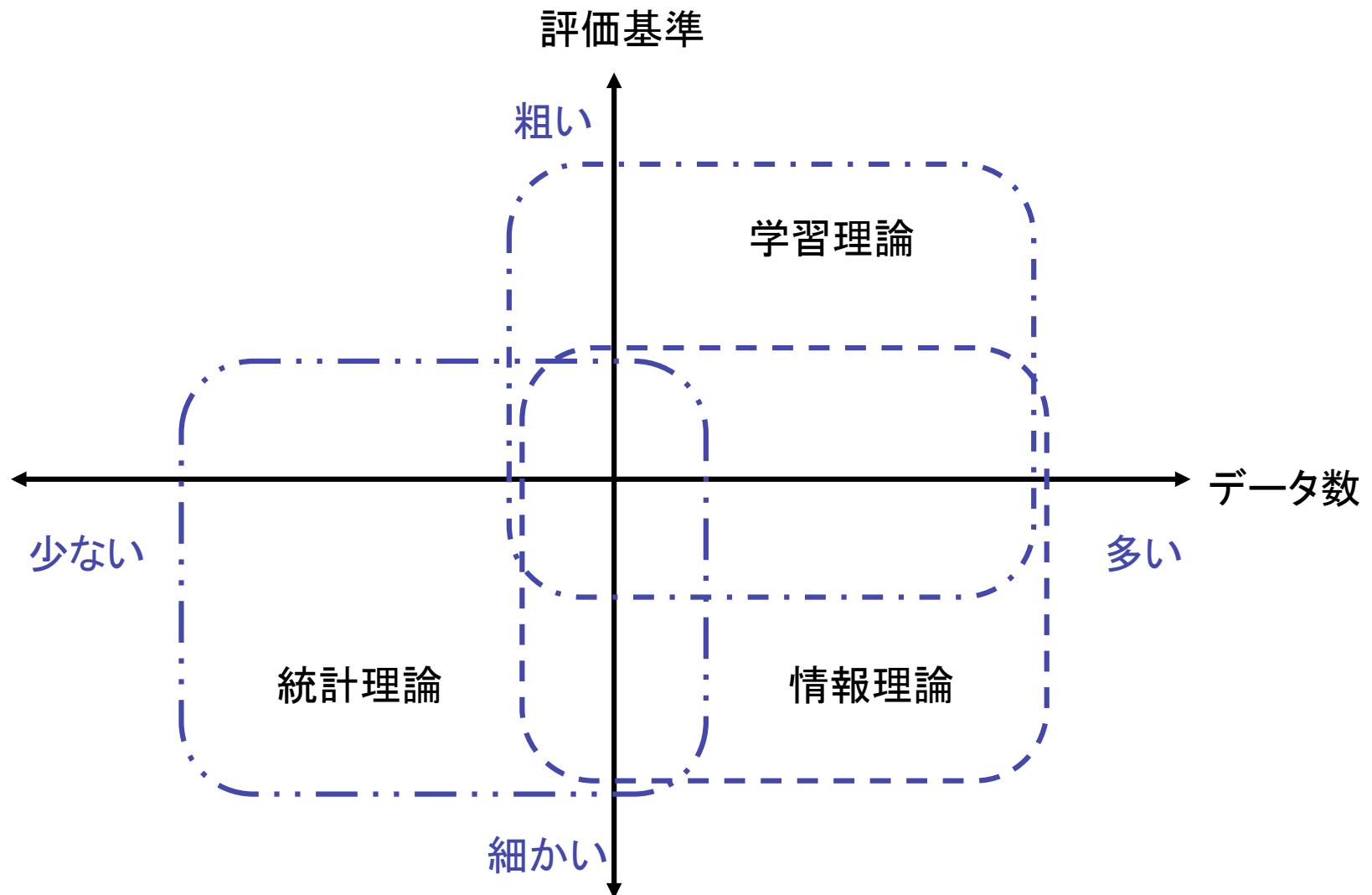
予測, 統計的複雑度,  
モデル選択, . . .

## 通信路符号化(冗長度を付加)

復号法(誤り訂正符号, CDMA, MIMO)

グラフィカルモデル,  
事後確率計算, BP,  
和積アルゴリズム, 平  
均場近似, . . .

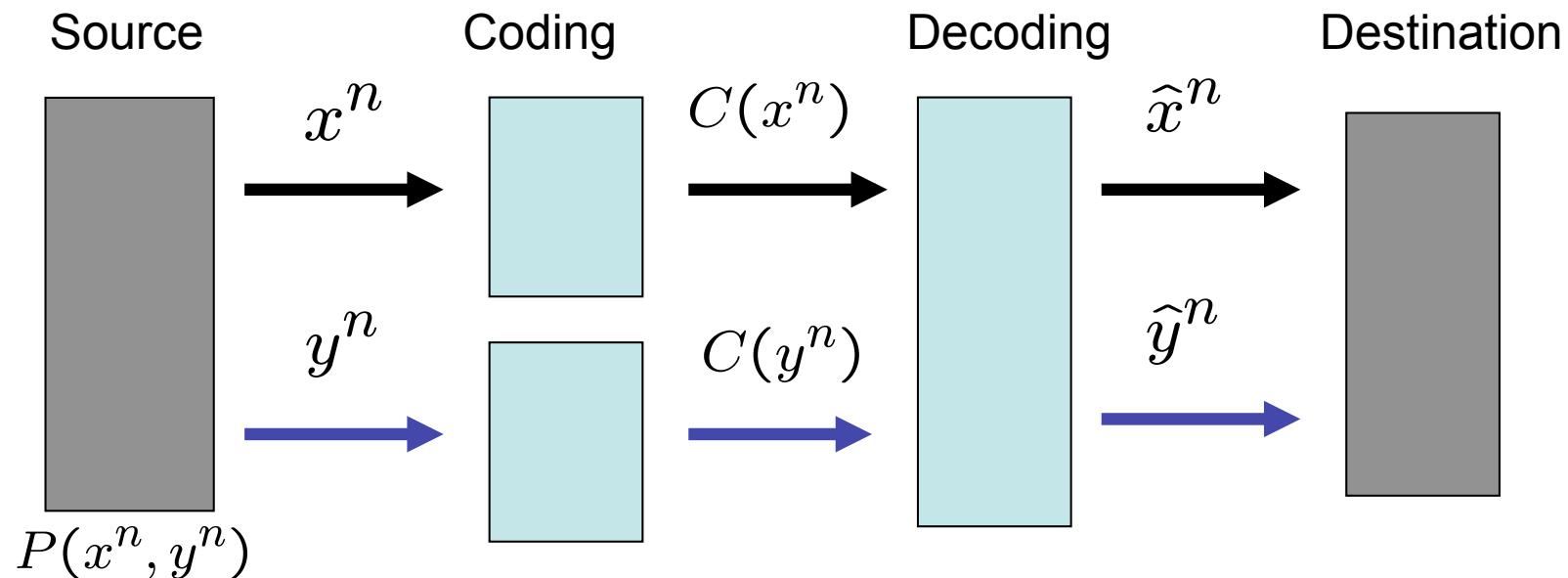
# 学習理論, 統計理論, 情報理論 (無理して特徴付けると)



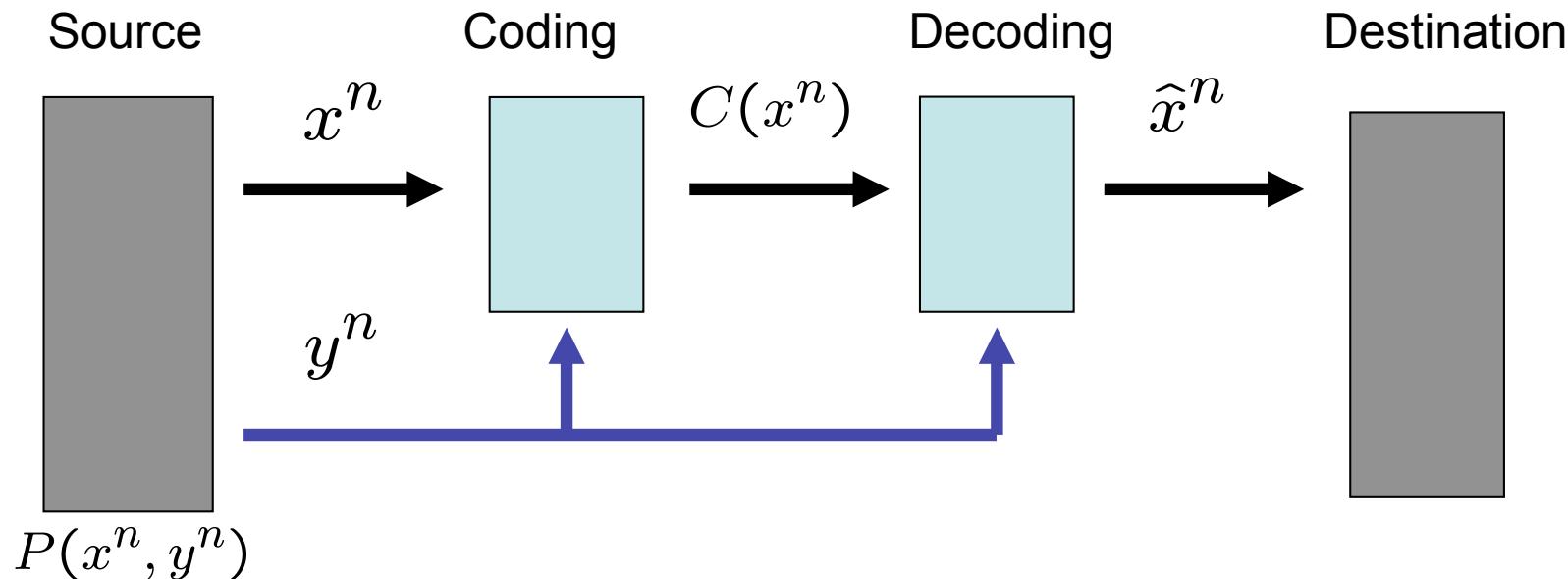
# 情報理論の特徴

- 評価基準が明確(復号誤り率, レート)
- 評価基準の限界をまず示す
- その限界を達成する符号化, 復号化を求める
- 実用的アルゴリズムを求める(計算量, メモリー量)

## 0. Slepian –Wolf Problem



# 1. Surece coding with side information



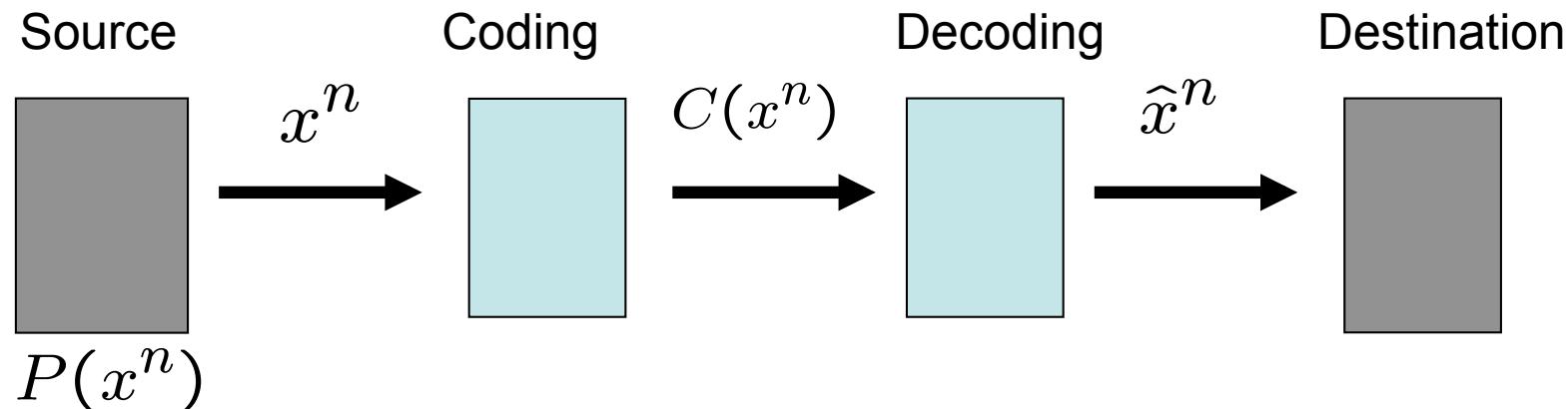
Distortion less:

$$x^n = \hat{x}^n$$

Mean code length:  $\lim_{n \rightarrow \infty} 1/nE[|C(X^n)|] \geq H(X|Y)$

## 2. Source coding

$$x^n = x_1, x_2, \dots, x_n (x_i \in A)$$



Distortion less:

$$x^n = \hat{x}^n$$

Mean code length:  $\lim_{n \rightarrow \infty} 1/nE[|C(X^n)|] \geq H(X)$

Ideal code length:

$$|C(x^n)| = -\log P(x^n)$$

If the distribution of a source is known, we can construct the optimal source code.

eg) Arithmetic codes

# ユニバーサル情報源符号化と学習理論

ブロック符号

$$C(x^n)$$

逐次符号

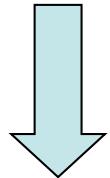
$$C(x_t|x^{t-1})$$

最適な符号化

$$|C(x^n)| = -\log P(x^n)$$

$$|C(x_t|x^{t-1})| = -\log P(x_t|x^{t-1})$$

ユニバーサル符号問題



$P(x_t|x^{t-1})$  が未知

損失関数:  $\log \hat{P}(x_t|x^{t-1}) - \log P(x_t|x^{t-1})$

累積損失: 符号長

累積リスク: 期待符号長

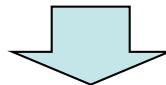
# 補助情報つき情報源符号化と学習理論

$$\begin{array}{c} x^n : x_1, x_2, \dots, x_n \\ | \quad | \quad | \\ y^n : y_1, y_2, \dots, y_n \end{array} \quad \begin{array}{|c|} \hline \text{目的変数} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{説明変数} \\ \hline \end{array} \quad (\text{補助情報})$$

$$\begin{array}{c} x_1, x_2, \dots, x_{t-1}, x_t \\ | \quad | \quad | \quad \uparrow \\ y_1, y_2, \dots, y_{t-1}, y_t \end{array} \quad \text{回帰問題, 分類問題}$$

$$\widehat{P}(x_t | y_t, x^{t-1}, y^{t-1})$$

# 補助情報つき情報源符号化の問題



情報理論的にはほとんど解決

補助情報つき情報源符号化についてはZiv[1984], Subarahmanya[1995], Muramatsu[1998], Yan[2001], Uematsu[2002]等の研究があり, i.i.d. や定常エルゴード情報源等について, 漸近最良の符号  $\frac{C \log n}{n}$  で平均符号長が条件付エントロピーに近づくことは示されている.

学習理論の問題を考えると, リスク関数を最小化するダイレクトで詳細な性能解析が望まれる(最速でエントロピーに近づく).

ベイズ符号で補助情報つき情報源符号化を考える

### 3. Why Bayes codes?

The advantage of Bayes codes?

1) Precise evaluation of the code length:

We can evaluate the exact asymptotic code length of Bayes codes as far as constant terms

2) The highest convergence rate:

The convergence rate of Bayes codes achieves the highest order for universal codes

3) Maximin codes and Minimax codes:

The class of Bayes codes includes Maximin codes (Jeffrys' prior) and Minimax codes

## 4-1. Bayes codes

Known  $\leftarrow$  a parametric distribution:  $P(x^n|\theta)$

Unknown  $\leftarrow$  k-dimensional real parameter vector:  $\theta \in \Theta \subset \mathbb{R}^k$

Under the condition that the **Arithmetic Coding** is used for coding, the decision of **encoding prob.**  $\hat{P}(x^n)$  is the main problem in universal coding.

The optimal solution from the view point of Bayes strategy

Block code:  $\hat{P}_B(x^n) = P(x^n) = \int_{\Theta} P(x^n|\theta)dP(\theta),$

A block sequence  $x^n$  is encoded to a code word  $C(x^n)$  simultaneously

Predictive code:  $\hat{P}_P(x_t|x^{t-1}) = P(x_t|x^{t-1}) = \int_{\Theta} P(x_t|x^{t-1}, \theta))dP(\theta|x^{t-1})$

A symbol  $x_t$  given  $x^{t-1}$  is encoded to a code word in order

## 4-2. the code length of Bayes codes

The code length of the Bayes **block** code is the same as that of the Bayes **predictive** code.

$$-\log \hat{P}_B P(x^n) = \sum_{t=1}^n -\log \hat{P}_P(x_t|x^{t-1})$$

Precise evaluation [Clarck&Barron90]:

$$E_\theta[-\log P(X^n)] = H(X^n|\theta) + \frac{k}{2} \log \frac{n}{2\pi e} + \log \frac{\sqrt{\det I_X(\theta)}}{p(\theta)} + o(1).$$

Fisher information

This convergence rate achieves the highest order for universal codes

## 5. Precise evaluation of the code length of the Bayes codes with side information

The average code length of Bayes codes with side information under some assumption is given by

$$\begin{aligned}
 & -E_{X^n Y^n} [\log P_c(x^n | y^n)] \\
 &= -E_{X^n Y^n} [\log P(x^n | y^n, \nu^*)] + \frac{k}{2} \log \frac{n}{2\pi e} \\
 &+ \frac{1}{2} \log \det I_{X|Y}(\nu^*) + \log \frac{1}{p(\nu^*)} + o(1),
 \end{aligned} \tag{1}$$

where  $f(\nu)$  is the prior probability of  $\nu$ .

Conditional Fisher information:

$$I_{X|Y}(\nu) = -\lim_{n \rightarrow \infty} \frac{1}{n} E_{X^n Y^n} \left[ \frac{\partial^2 \log P(X^n | Y^n, \nu)}{\partial \nu (\partial \nu)^T} \right]$$

$$E_\theta [-\log P(X^n)] = H(X^n | \theta) + \frac{k}{2} \log \frac{n}{2\pi e} + \log \frac{\sqrt{\det I_X(\theta)}}{p(\theta)} + o(1).$$

## 6. Bayes codes (unknown model )

Known  $\leftarrow$  a parametric distribution class:  $P(x^n|\theta(m), m)$

Unknown  $\leftarrow$  k-dimensional real parameter vector:  $\theta \in \Theta \subset \mathbb{R}^k$

Unknown  $\leftarrow$  parametric model:  $m \in M$

The encoding probability of Bayes codes:

$$\hat{P}_B(x^n) = \sum_{m \in M} \int P(x^n|\theta_m, m)p(\theta_m|m)P(m)d\theta.$$

Expectation is taken all over the model class

Precise evaluation for the code length for unknown model problem is given by[Gotoh & Matsushima98]

## The disadvantage of Bayes codes?

The complexity of the coding calculation

Assuming **the conjugate prior** as the prior probability of  $\theta$ ,  
the integral does not need for calculating the coding probability.

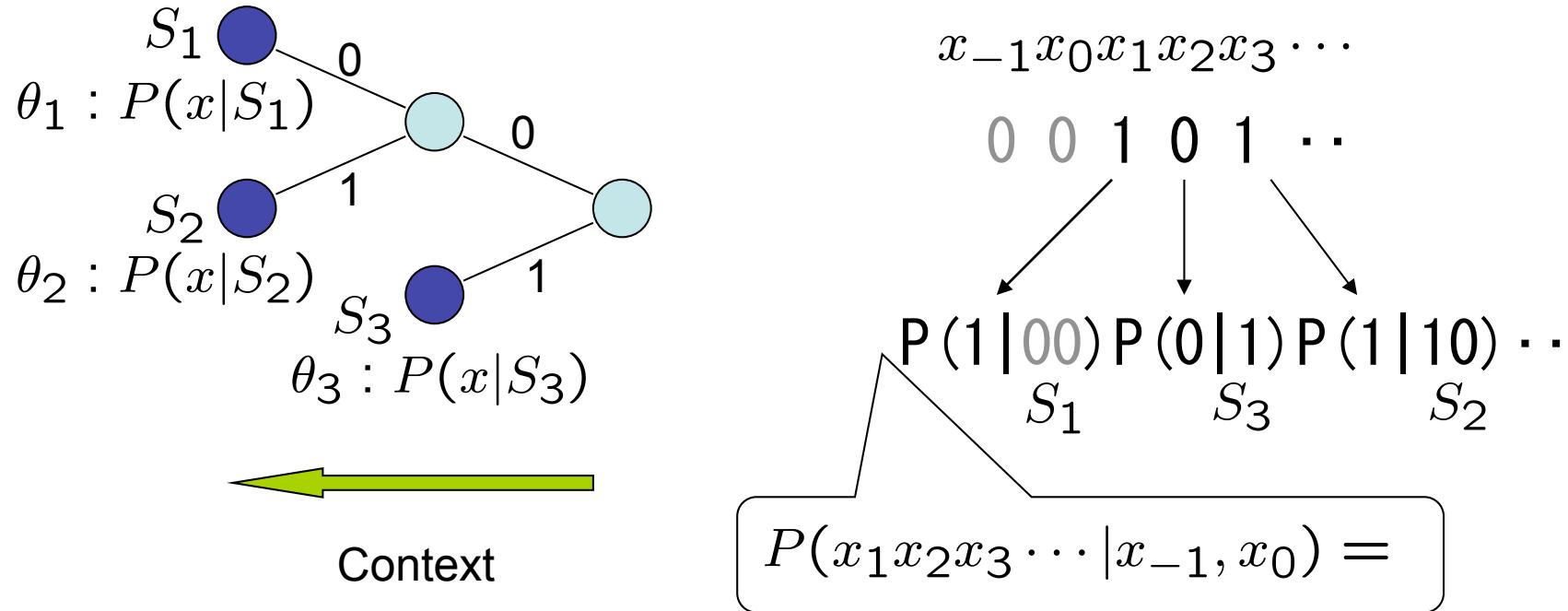
Integral is taken all over the parameter space

$$\hat{P}_B(x^n) = \sum_{m \in M} \int_{\Theta} P(x^n | \theta_m, m) p(\theta_m | m) P(m) d\theta.$$

Summation is taken all over the model class

Assuming **a special prior distribution** on the models, some efficient algorithms can be constructed.

## 7-1. Context tree models



A context tree(FSMX) model  $m$  is represented by a context tree( or a set of contexts) and a set of the parameters on the leaves :

$$P(x^n|m, \theta(m))$$

$\swarrow$        $\searrow$

$\{S_1, S_2, S_3\}$

$\{\theta_1, \theta_2, \theta_3\}$

## 8. Bayes codes for Context tree models(1)

Unknown model  $m_i = \{S_1, S_2, S_3 \dots\} \in M$

Unknown parameters  $\theta(m) = \{\theta_1, \theta_2, \theta_3 \dots\} \in \Theta(m)$

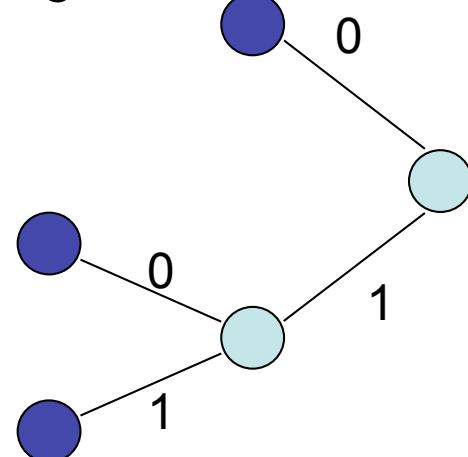
$M$ : the class of context trees up to depth 2

$m_2$ : simple Markov model

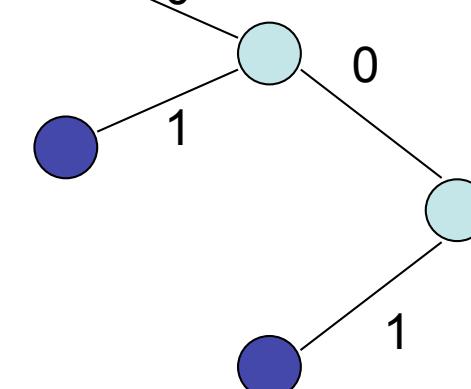
$m_1$ : i.i.d



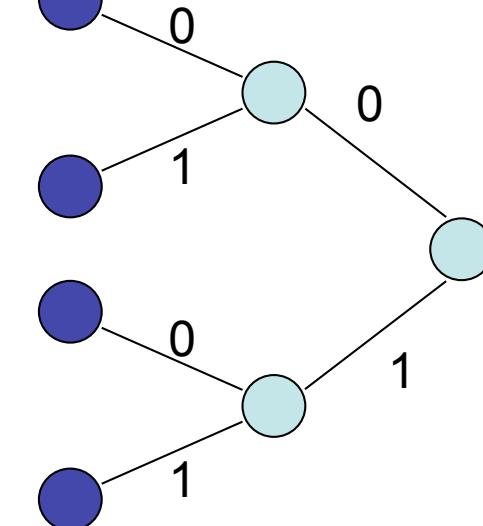
$m_3$



$m_4$



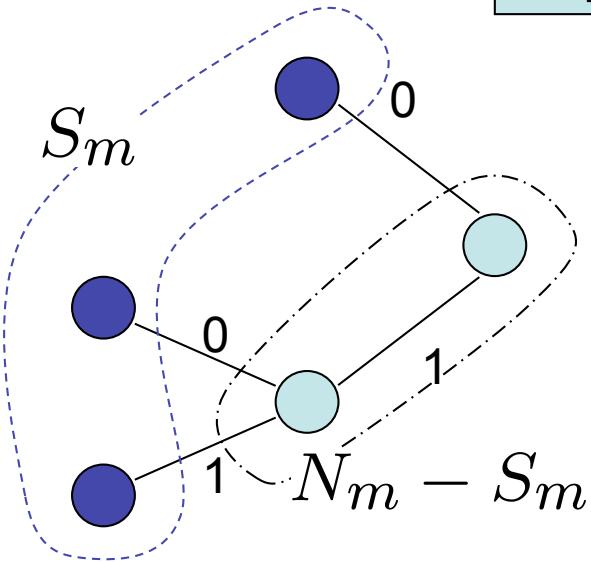
$m_5$



For calculation of the Bayes code, the expectation is taken over all context tree models whose number is  $O(2^{|X|^{d-1}})$

## 8. Bayes codes for Context tree models(2)

A special class of the prior distribution of  $m$



We assume that the prior probability  $P(m)$  is represented by the following formula:

$$P(m) = \prod_{s \in N_m - S_m} (1 - q(s)) \prod_{s_L \in S_m} q(s_L). \quad (1)$$

where

$$q(s) = \frac{P(S_{m_s})}{\sum_{S \in S_d} P(S_{m_s} \cup S)}, \quad (2)$$

where  $m_s \in \{m | s \in S_m\}$ ,  $S_d$  denotes the class of the set of descendant nodes of  $s$  that construct a complete tree.

The prior probability of  $m$  is represented by  $\{q(s) | s \in N_m\}$ . Moreover, the posterior probability  $P(m)$  is also represented by  $\{q(s|x^t) | s \in N_m\}$ .

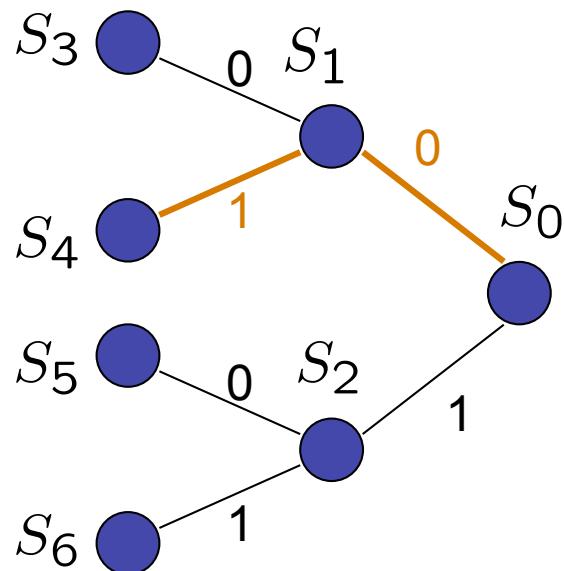
The CTW algorithm is regarded as the Bayes code assuming  $q(s) = 1/2$

## 8. Bayes codes for Context tree models(3)

The efficient coding algorithm is given by using the tree to which all context tree models in  $M$  are merged and the special prior.

$$\hat{P}(x^n) \leftarrow \text{Block codes: CTW algorithm}$$

$$\hat{P}(x_{t+1}|x^t) \leftarrow \text{Predictive codes: Predictive Bayes algorithm for Context tree models [Matsushima95]}$$



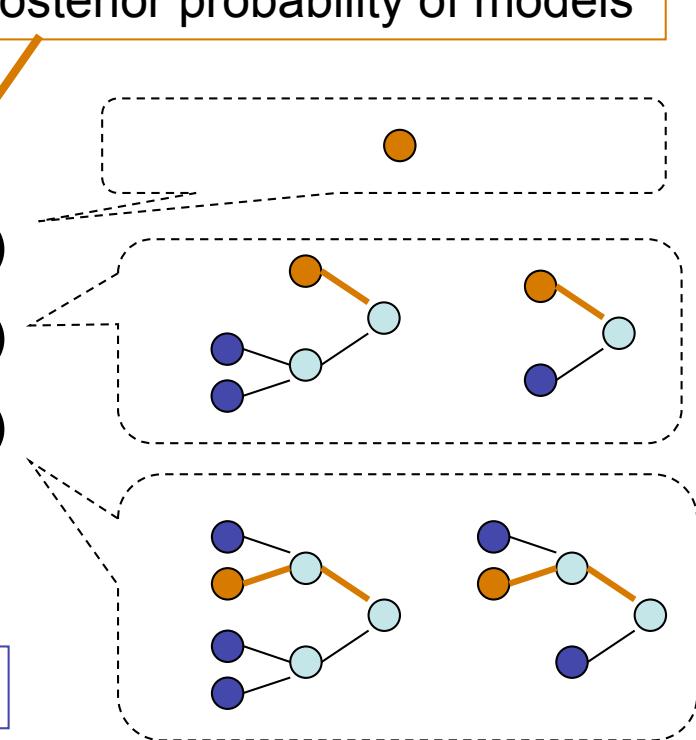
$\cdots x_{t-1}x_tx_{t+1}$   
 $\dots 1 \quad 0 \quad 0$

$$\begin{aligned} P(0|10) &= \hat{P}(0|S_0)\hat{q}(S_0|x^t) \\ &+ \hat{P}(0|S_1)\hat{q}(S_1|x^t) \\ &+ \hat{P}(0|S_4)\hat{q}(S_4|x^t) \end{aligned}$$

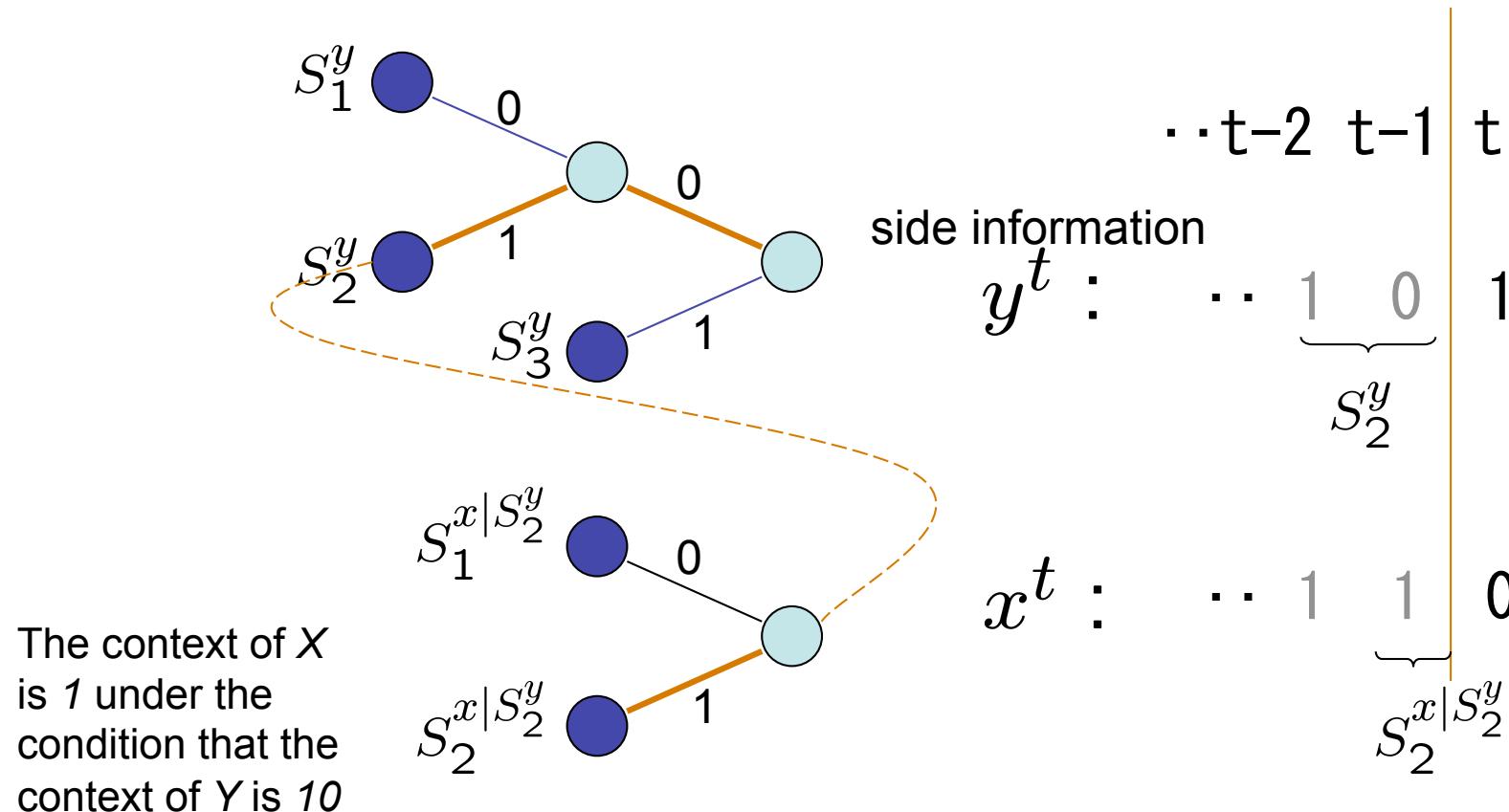
$$\hat{P}(0|S_i) = \frac{n(0|S_i) + \alpha/2}{n(S_i) + \alpha}$$

Posterior probability of models

By using  $q(s)$ , the complexity is deduced to  $O(d)$ .



## 9. Side information Context tree models



A side information context tree model is represented by

$$P(x_t|x^{t-1}, y^{t-1}) = P(x_t|S^{x|S^y}, \nu)P(S^y|\mu)$$

## 9. A special prior on Side-information context tree models

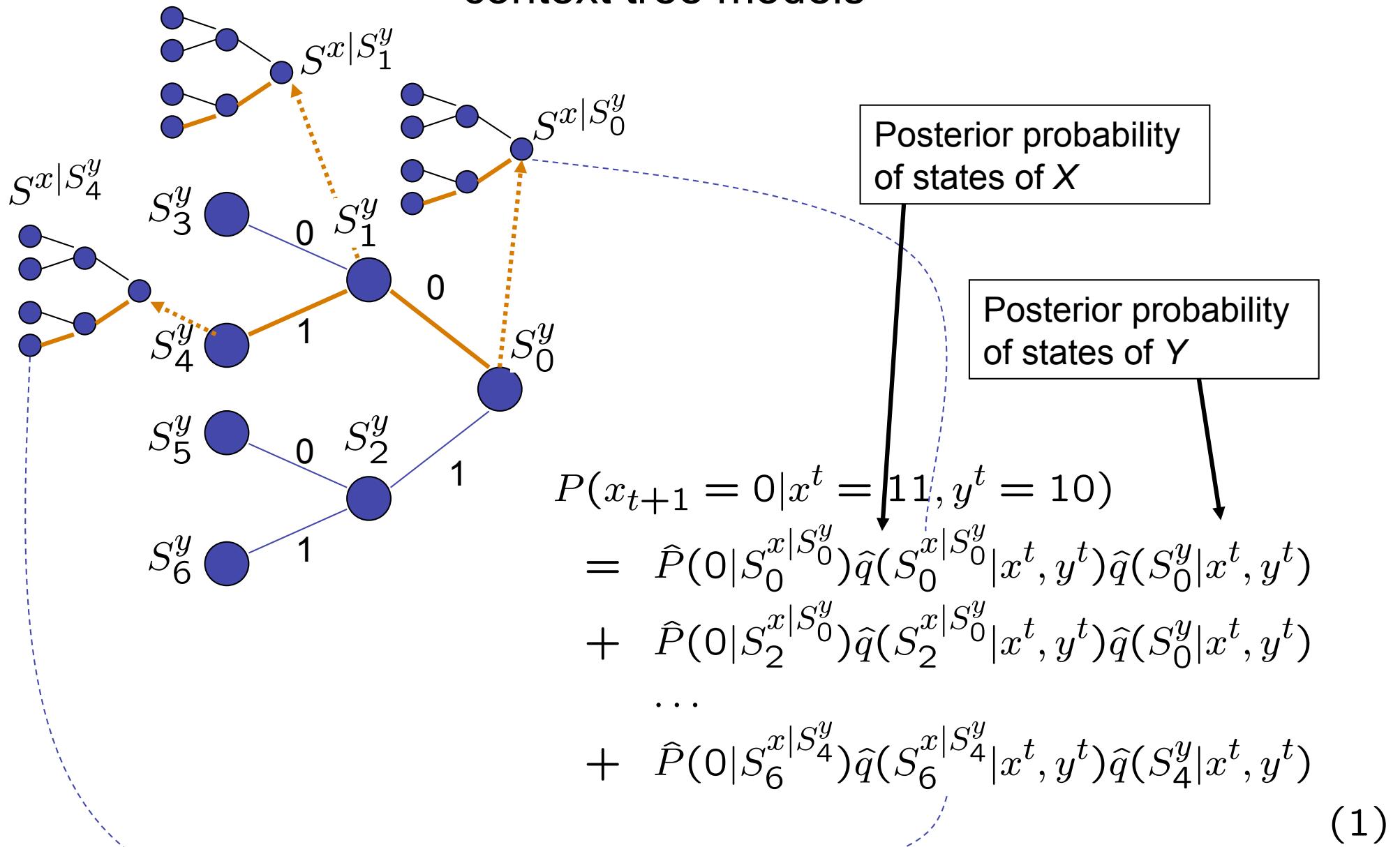
A special class of the prior distribution of  $m$

[Matsushima 05]

We assume that the prior probability  $P(m)$  of a side information context tree model is represented by the following formula:

$$\begin{aligned} P(m) &= \prod_{s^y \in N_m^y - S_m^y} (1 - q(s^y)) \prod_{s_L \in S_m^y} q(s_L^y) \\ &\quad \prod_{s^x \in N_m^{x|s_L^y} - S_m^{x|s_L^y}} (1 - q(s^x | s_L^y)) \prod_{s_L^x \in S_m^{x|s_L^y}} q(s_L^x | s_L^y). \end{aligned} \tag{1}$$

## 10. Algorithm of Bayes codes for Side-information context tree models



# 11. Precise evaluation of the code length of the Bayes codes

The average code length of Bayes codes with side information under some assumption is given by

$$\begin{aligned}
 & -E_{X^n Y^n} [\log P_c(x^n | y^n)] \\
 &= -E_{X^n Y^n} [\log P(x^n | y^n, \nu^*)] + \frac{k}{2} \log \frac{n}{2\pi e} \\
 &+ \frac{1}{2} \log \det I_{X|Y}(\nu^*) + \log \frac{1}{p(\nu^*)} + o(1),
 \end{aligned} \tag{1}$$

where  $f(\nu)$  is the prior probability of  $\nu$ .

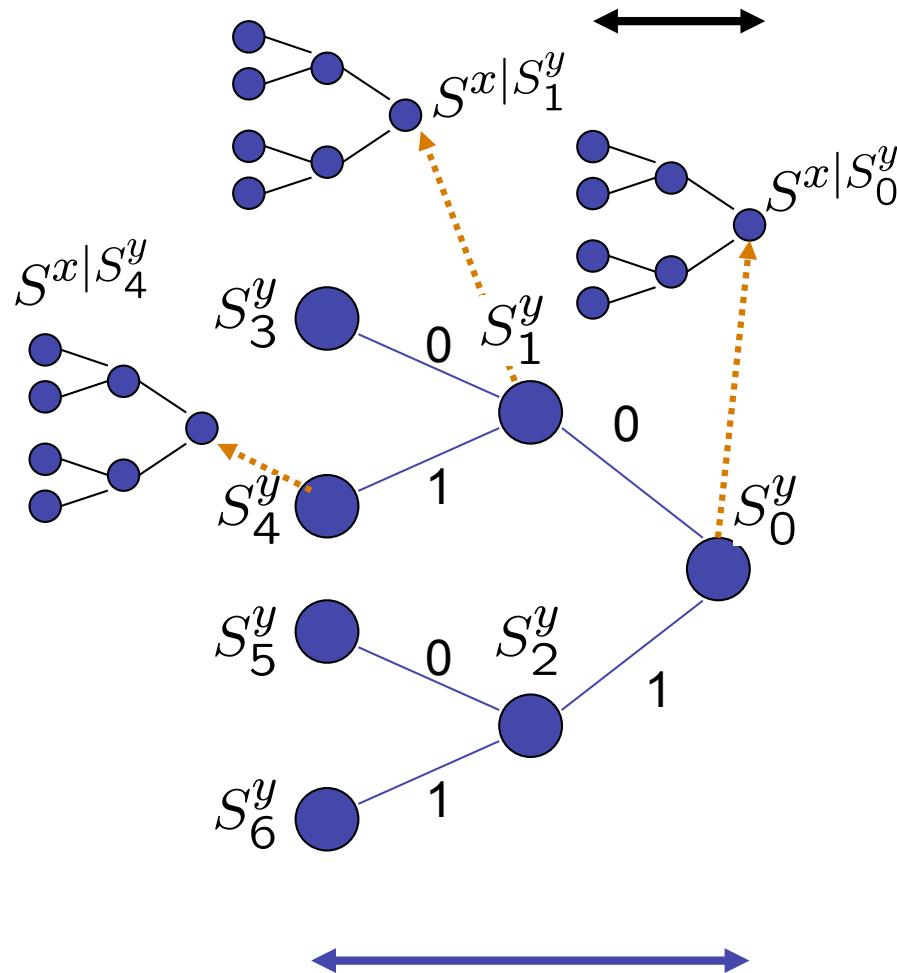
Conditional Fisher information:

$$I_{X|Y}(\nu) = -\lim_{n \rightarrow \infty} \frac{1}{n} E_{X^n Y^n} \left[ \frac{\partial^2 \log P(X^n | Y^n, \nu)}{\partial \nu (\partial \nu)^T} \right]$$

$$E_\theta [-\log P(X^n)] = H(X^n | \theta) + \frac{k}{2} \log \frac{n}{2\pi e} + \log \frac{\sqrt{\det I_X(\theta)}}{p(\theta)} + o(1).$$

## 12. Time complexity of the algorithm

$d(x)$  : The depth of the context tree for  $x$

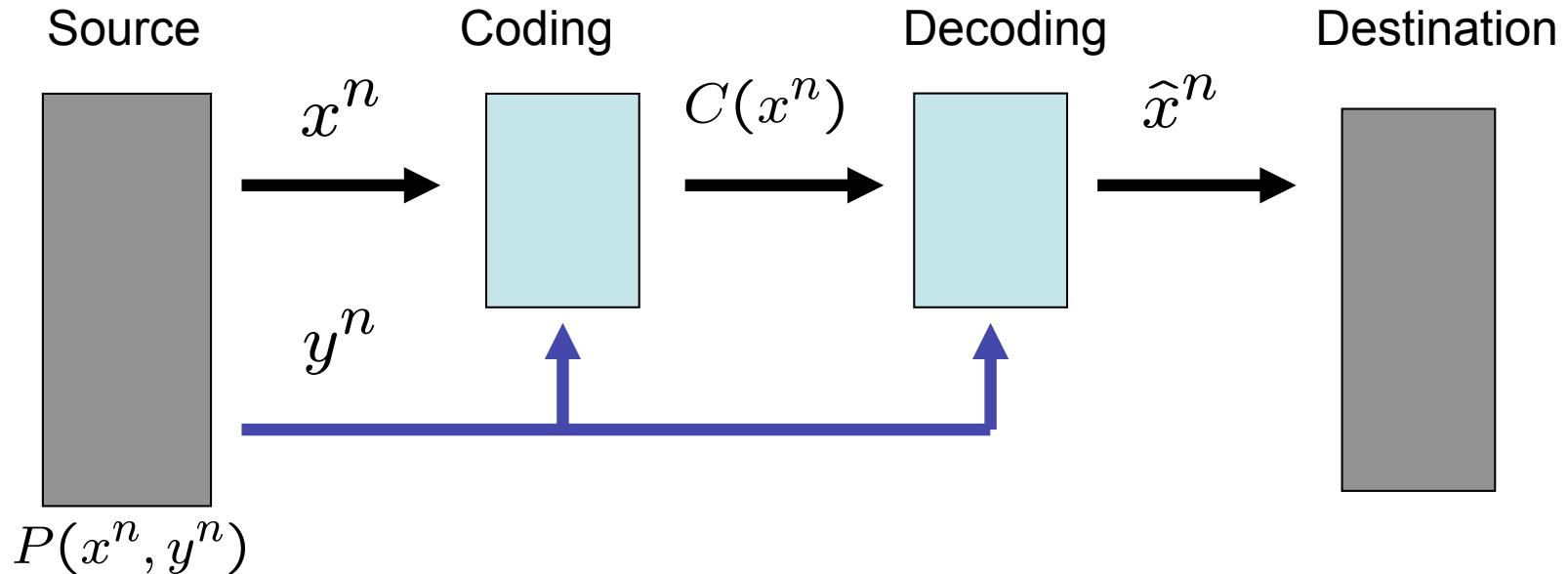


The time complexity of the algorithm:  
 $O(d(x)d(y))$

The parallel time complexity of the  
algorithm:  $O(d(x)+d(y))$

$d(y)$  : The depth of the context tree for  $y$

# まとめ



- ・補助情報付ユニバーサル符号と学習理論の関係
- ・補助情報付ベイズ符号の累積log損失(符号長)の漸近評価
- ・補助情報文脈木モデルとその上での効率的符号化アルゴリズム